

Calibrating Recurrent Sliding Window Classifiers for Sequential Supervised Learning

Saket Joshi & Thomas G. Dietterich
School of Electrical Engineering and Computer Science
Oregon State University
Corvallis, OR 97331

November 21, 2003

Abstract

Sequential supervised learning problems involve assigning a class label to each item in a sequence. Examples include part-of-speech tagging and text-to-speech mapping. A very general-purpose strategy for solving such problems is to construct a recurrent sliding window (RSW) classifier, which maps some window of the input sequence plus some number of previously-predicted items into a prediction for the next item in the sequence. This paper describes a general-purpose implementation of RSW classifiers and discusses the highly practical issue of how to choose the size of the input window and the number of previous predictions to incorporate. Experiments on two real-world domains show that the optimal choices vary from one learning algorithm to another. They also depend on the evaluation criterion (number of correctly-predicted items versus number of correctly-predicted whole sequences). We conclude that window sizes must be chosen by cross-validation. The results have implications for the choice of window sizes for other models including hidden Markov models and conditional random fields.

1 Introduction

1.1 Sequential Supervised learning

The standard supervised learning problem is to learn to map from an input feature vector x to an output class variable y given N training examples of the form $(x_i, y_i)_{i=1}^N$. Problems where the variable y can take up a finite number of discrete values are called classification problems in the literature. Standard supervised learning algorithms like decision tree algorithms, the naive Bayes algorithm, the K-nearest neighbor algorithm etc., learn this map and come up with a classifier which given a new example feature vector x_{new} as input, produces the corresponding class value y as output. The accuracy of a classifier is measured as the percentage of such new data points (data points not included in the set of training examples) it classifies correctly.

Many problems, however, do not succumb to this simple approach. Consider the problem of part-of-speech tagging. The problem can be defined as follows. Given a set of sentences in the English language, label each word according to the part of speech it represents in the sentence. Trying the classical supervised learning approach on this problem would lead to construction of a dataset where each example would be a word broken down into a number of features and the class variable specifying the corresponding part of speech (e.g., noun, verb, adjective, etc.). We could build a classifier by feeding some training data (data where we already know the correct part of speech for each word) to a standard learning algorithm and then classify new words. This method would not work well however, as the problem arises with words like “lead”. There is no way to tell if the word “lead” viewed in isolation is an adjective as in “a lead pencil”, or a verb as in the sentence “please lead the way”. The classification here depends on the relationship of this word to the other words in the sentence. The classification is bound to the “sequence” in which this word appears. Therefore the part-of-speech can be seen as a problem of mapping input sequences to output sequences. No learning algorithm can be accurate in this setting unless it has a way to analyze the context around individual words. A standard Machine Learning approach would try to process

each word separately, but this example shows that would not work because at a minimum each word needs to be handled in context. We need algorithms to learn a map from input sequences to output sequences.

Many recent learning problems can be viewed as extensions of standard supervised learning to the setting where each input object X is a *sequence* of feature vectors, $X = \langle x_1, x_2, \dots, x_T \rangle$, and the corresponding output object Y is a sequence of class labels $Y = \langle y_1, y_2, \dots, y_T \rangle$. The *sequential supervised learning* (SSL) problem is to learn to map from X to Y given a set of N training examples, $(X_i, Y_i)_{i=1}^N$.

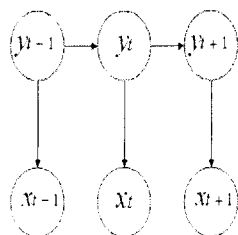


Figure 1: Hidden Markov Model

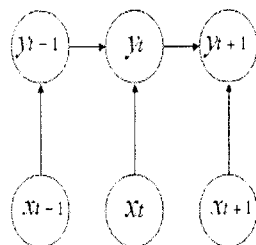


Figure 2: Max Entropy Markov Model

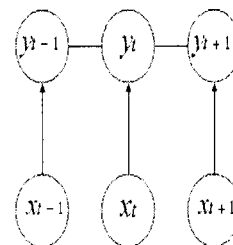


Figure 3: Conditional Random Fields

1.2 Popular sequential supervised learning algorithms

In the literature, two general strategies for solving SSL problems have been studied. One strategy, which we might call the direct approach, is to develop probabilistic models of sequential data. The advantage of these probabilistic models is that they seek to capture the true sequential relationships that generate the data. Examples of such strategies are hidden Markov models, maximum entropy Markov models, and conditional random fields. The other general strategy that has been explored might be called the Indirect approach (i.e., a “hack”). In this strategy, the sequential supervised learning problem is solved indirectly by first converting it into a standard supervised learning problem, solving that problem, and then converting the results into a solution to the SSL problem. Examples of indirect approaches include sliding windows and recurrent sliding windows. Let us look at each of them with more emphasis on sliding windows and recurrent sliding windows because they are the basis for this paper.

1.2.1 Hidden Markov Models

The hidden Markov model (Rabiner, 1989) is a generative model of the joint distribution $P(X, Y)$ of the object sequence X and the label sequence Y . Dynamic Bayesian networks can use information from previous time steps (or neighboring examples in the sequence) for current evaluation. Dynamic Bayesian networks factor the state space into subspaces. This is especially useful when the interactions between the subspaces is sparse. If the subspaces are highly correlated, it is more useful to combine them and treat them as a single space. A hidden Markov model is a special case of dynamic Bayesian networks that does not factor the state space. Therefore it is used in problems where such factoring is superfluous. Figure 1 shows an example of a HMM. The popular forward-backward prediction algorithm provides efficient inference for HMMs. We will see examples of real world problems currently solved using HMMs in the next chapter.

1.2.2 Maximum Entropy Markov Models

The HMM is a generative model, which means that it models the relationship (c) as a map from y to x . During classification, Bayes’ rule is employed to infer the y values given the x values. There are two potential problems with this method. Firstly a joint distribution of the object sequence and the label sequence has to be maintained because of which, the required resources grow exponentially with the number of input features. Secondly a generative model may not exactly model the causal relationships between the x ’s and the y ’s correctly. In some problems the causal relationship is from the x ’s to the y ’s. Maximum entropy Markov models (McCallum, Freitag & Pereira, 2000), which are another sequential supervised learning device, model exactly that relationship.

Table 1: Simple sliding windows

(X, Y)	enough	In ^{^-} -f-
w_1	___enou	I
w_2	__enoug	n
w_3	_enough	^
w_4	enough_	-
w_5	nough__	f
w_6	ough___	-

Table 2: Recurrent sliding windows

(X, Y)	enough	In ^{^-} -f-
w_1	___enou__	I
w_2	__enoug_I	n
w_3	_enoughIn	^
w_4	enough_n^	-
w_5	nough__^-	f
w_6	ough___-f	-

MEMMs learn conditional probability tables rather than a joint distribution. This avoids problems caused by a large input feature space. MEMMs require only simple extensions to the forward-backward algorithm. Figure 2 shows an example of a MEMM.

1.2.3 Conditional Random Fields

More recently, the conditional random field (Lafferty, McCallum, & Pereira, 2001) has been proposed as a model of the conditional distribution $P(Y|X)$. It is another sequential supervised learning approach that goes beyond MEMMs in that it works with a conditional model and additionally solves the label bias problem that is caused by the use of the exponential transition model of MEMMs. In MEMMs, the probabilities of state transitions are normalized locally for each time step rather than normalizing over the entire sequence of y 's. As a result, they do not compute the globally most likely sequence. This means that the probabilities of transitions from state y_t to state y_{t+1} are normalized over all possible states y_{t+1} that can be reached from y_t . As a result, probabilities of transitions to states with fewer outgoing states increase. This gives rise to a bias towards states with fewer number of possible next states, and this is termed as the label bias problem. CRFs solve this problem by employing global normalization. CRFs can be pictured as Markov random fields of the y 's conditioned on the x 's. An example of CRFs is shown in Figure 3. Conditional random fields constitute a very promising approach, and they are attracting further applications and research.

1.2.4 Sliding Windows and Recurrent Sliding Windows

The indirect approach employed in sliding windows and recurrent sliding windows is to convert the input and output sequences into a set of windows as shown in Table 1. Each window consists of central element x_i and *LIC* letters of left input context and *RIC* letters of right input context (in the figure $LIC = RIC = 3$). Contextual positions before the start of the sequence or after the end of the sequence are filled by a designated null value (in this case _).

In most SSL problems, there are regularities in the sequence of y values. In part-of-speech tagging, the grammar of natural language constrains the possible sequences of parts of speech. In text-to-speech mapping, there are patterns in the phoneme sequence. The simple sliding window method cannot capture these patterns unless they are completely manifested in the X values as well, which is rarely the case. One way to partially learn these patterns is to employ *recurrent* sliding windows, in which previous predictions (e.g., for y_{t-1}, y_{t-2} , etc.) are fed back as input features to help predict y_t . During learning, the observed labels in the training set can be used in place of these fed back values. During classification, the sequence is processed from left-to-right, and the predicted outputs $\hat{y}_{t-2}, \hat{y}_{t-1}$ are fed as input features to predict y_t . Note that the sequence could also be processed from right-to-left, but not simultaneously left-to-right and right-to-left. We denote the number of fed back y values as the left output context *LOC* or right output context *ROC*, depending on the direction of processing. Table 2 shows the training windows with $LOC = 2$ and $ROC = 0$. The advantage of the recurrent sliding window approach is that it can be combined with *any* standard supervised learning algorithm to solve SSL problems.

1.3 The 4-Case Analysis for sequential supervised learning problems

It is hard to predict an output sequence given an input sequence directly in a single step. The number of possible output sequences is exponential in the length of the sequences, and classification methods do not work well with

Table 3: The 4-Case Analysis							
		RSW		HMM		CRF	
b	d	IC	OC	IC	OC	IC	OC
0	0	0	0	0	0	0	0
0	1	>0	0	0	>0	>0	0
1	0	>0	>0	0	>0	0	>0
1	1	>0	>0	0	>0	>0	>0

large numbers of classes. Therefore we seek a divide-and-conquer approach. There are four relationships that are important in a sequential supervised learning problem:

- (a) The sequential relationships among the x 's
- (b) The sequential relationships among the y 's
- (c) The mapping from x 's to y 's
- (d) The relationship between the y_t 's and the x_{t-1} 's or the x_{t+1} 's

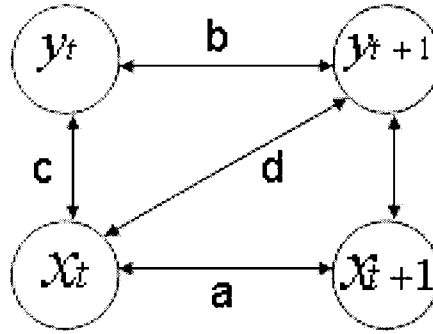


Figure 4: Pictorial View of the 4 Relationships in Data

In a given dataset, the x 's are observed values. Therefore unless some values are missing, and we are interested in filling in the blanks, relationship (a) is irrelevant. Relationship (c) is essential for a problem to qualify as a supervised learning problem. In the absence of relationship (c) the data represents a time series or Markov chain problem. Hence we assume that relationship (c) is present. Based on the presence or absence of relationships (b) and (d), we can imagine 4 cases as shown in Table 3. Table 3 also shows expected sizes of input and output contexts for different learning algorithms on problems that fall in each of these four categories. RSW encompasses recurrent sliding windows and sliding windows. HMM represents hidden Markov models and other generative sequential data models. CRF represents conditional random fields, maximum entropy Markov models and other conditional sequential data models. In the 00 case, the data is non-sequential and hence we can expect neither input nor output context to be useful in any of the methods. In the 01 case, $x_{t-n} \dots x_{t+n}$ contain information about y_t . This relationship can be directly captured by an input context in RSW and CRF. Hence we expect no output context there (because relationship (b) is absent). HMMs however by virtue of being generative models, cannot have an input context. Therefore we expect them to employ an output context and try to approximate relationship (d) as the product of relationship (c) and relationship (b). In the 10 case, relationship (d) is absent and relationship (b) can be exploited by an output context alone. We expect such behavior from HMM and CRF. RSW might employ an input context and try to approximate relationship (b) as the product of relationship (c) and relationship (d). Note that relationships (b), (c), and (d) form a triangle and any method that learns two

of these relationships learns an approximation of the third. In the 11 case, we expect all learning algorithms to employ both input and output contexts (Except HMMs for which there is no input context).

1.4 Our Approach

Two practical issues arise in applying either the direct or indirect methods. The first issue is the size of the context. How much of the input sequence and output sequence should be considered at each point along the sequence? For the recurrent sliding window methods, this can be stated concretely as follows: how large should the left input context, right input context, and left output context be? Presumably larger contexts will lead to high variance and overfitting by the learning algorithm, whereas smaller contexts will lead to lower variance but higher bias and underfitting by the learning algorithm.

Based on the 4-case analysis, we can say that only in cases 10 and 11 do we expect there to be a trade-off between input and output contexts. In the 01 case, output context should always be 0. So in this case there is no trade-off: output context is simply useless. This may explain some of our results (or at least it may suggest a way of interpreting our results).

The second issue is the criteria for measuring performance. A straightforward generalization of the 0/1 loss for standard supervised learning is to measure the number of whole sequences that are correctly predicted. However, another measure is the number of individual x_t items correctly classified (e.g. individual words in part-of-speech tagging). An interesting question is whether the optimal choice of input and output context depends on the performance criterion. Let us analyze what the 4-case analysis says about this.

- 00: The two criteria should be the same.
- 01: If we use a very large input context, this will make the whole word predictions more accurate, but it may also increase variance. So we might predict that a smaller input context would be more appropriate for individual letter predictions, since whole word accuracy is less important and this pushes us to reduce variance by using a smaller input context.
- 10: Again we predict that for whole word accuracy we want larger output context and a smaller context would be appropriate for single letter predictions.
- 11: For the whole word criterion, this would suggest larger output context (which, because of the trade-off means smaller input context). For the single letter criterion it depends on which influence is stronger. If the $y \rightarrow y$ correlations are weaker than the $x \rightarrow y$ correlations, then we use a larger input context, otherwise a larger output context.

In order to test our predictions we conducted an experimental study of these two issues. In this paper, we present the experimental study employing a general-purpose recurrent sliding window system, RSW, that we have constructed as part of WEKA, the widely used Java based Machine Learning library developed at the University of Waikato, NZ. RSW can work with any of the WEKA classifiers that return class probability distributions. Our study compares naive Bayes, decision trees, bagged trees, and adaboosted trees on two SSL problems: NETtalk and protein secondary structure prediction. We find that our intuitions are somewhat justified. There is some trade-off between the amount of input context and the amount of output context for the NETtalk data (which clearly belongs to case 11), but it is not crystal clear. Furthermore, we show that better learning algorithms (e.g., adaboosted trees) are able to handle larger amounts of both input and output context. Finally, we show that the appropriate choice of input and output context depends on the learning algorithm and on the evaluation criterion. This suggests that there are no general-purpose methods for choosing the amount of input and output context in an algorithm-independent way.

2 Previous Work

At this point we will present an introduction of some previously studied sequential supervised learning problems and approaches that have been used to solve them.

2.1 NETtalk (text-to-speech mapping)

The NETtalk problem is defined as follows. Given a list of words and their corresponding phoneme-stress sequences, learn a map that can be applied to predict the correct phoneme-stress sequence of a new word. A phoneme is the most basic unit of speech and a stress is a measure of emphasis given to that phoneme during pronunciation. The phoneme and stress pair together form a label which can be fed to a speech synthesis system to generate spoken words. Since the words themselves can be viewed as a sequence of letters, NETtalk can be viewed as a sequential supervised learning problem in which the x 's are the individual letters and the y 's are the corresponding phoneme-stress pairs (class labels). Some of the early work with sequential data using sliding windows was done for the NETtalk system (Sejnowski & Rosenberg, 1987). The NETtalk system learnt to pronounce English text by learning a map from individual letters in a word to their respective phonemic sound symbols coupled with the amount of stress associated with that phoneme. An example is given in Tables 1 and 2 in the previous section. Sejnowski & Rosenberg (1987) used a neural network with 203 input units, 80 hidden units, and 26 output units. The input units were divided into sets of 29 units each (26 letters and 3 extra features). Each feature was a single bit that was 1 if a particular letter appeared at a particular position in the input. Each of these 7 sets represented a letter and altogether a 7 letter window. The NETtalk system achieved a performance of around 93% on the training data for predicting individual letters correctly and around 65% for predicting whole words correctly. However this system was specifically built and tuned to solve the text-to-speech problem alone. An interesting inference with the NETtalk data that was indicated by Sejnowski & Rosenberg, and later clearly presented by Bakiri & Dietterich (1991), was that the pronunciation (phoneme and stress) of a letter depends more on the letters that come after it than on those that come before it. As a consequence, a window sliding over the word from right to left gives better performance than a window sliding in the opposite direction. Later Adamson and Dampier (1996), employed a similar approach but with a recurrent network on this problem of learning to pronounce English text. We have used the NETtalk data for our experiments.

2.2 Prediction of protein secondary structure

An important problem in molecular biology is that of predicting "protein secondary structure". Briefly, the task is as follows. For each element in a protein sequence (drawn from a 20-letter alphabet), assign it to one of three classes (alpha helix, beta sheet, or coil). If we view the class labels or residues (alpha helix, beta sheet, or coil) as a sequence corresponding to the protein sequence, then this task can be formulated as a sequential supervised learning problem where the x 's represent the proteins and the y 's represent the individual class label (alpha helix, beta sheet, or coil). As a result of the sequential nature of protein and DNA data, applications of sequential data classifiers in Computational Biology are immense. Qian & Sejnowski (1988) employed a sliding window of width 15 with neural networks in their work on prediction of secondary structure of globular proteins. Later Krogh et al. (1993) used an HMM to model protein families. Again the attempts here have been to come up with a model specifically tuned to solve this problem. We have used the Qian & Sejnowski data sets for training and testing of prediction of secondary structure of proteins as well, in our experiments.

2.3 Part-of-speech tagging

Part-of-speech tagging has been defined in the previous section. For more than thirty years POS tagging has been a standard for the sequential data learning task. Yet solutions to this problem have been ad-hoc. Words that can potentially belong to more than one part of speech category are called ambiguous words in the literature. Researchers have figured out to a large extent what the unambiguous words are and have built dictionaries for them. The majority of the research in POS tagging is conducted for disambiguation given the context in which the word appears. The context is generally selected as a few words appearing before and after the current word in the sentence. Greene & Rubin (1971) used a rule-based approach to build a text tagger and achieved up to 77% accuracy on test data. Later improvements were made, and statistical methods started finding their way into this area quite early. Jelinek (1985) used a HMM for this purpose. Almost all work on POS tagging after that has been with improving HMM performance for this problem. Cutting, Kupiec et al. (1992) also used a HMM for their POS tagging system and employed ambiguity classes (classes of ambiguous words) and a phrase recognition system based on word sequences to achieve an accuracy of 96% on unseen data. Brill (1992) deviated from the main stream and introduced a trainable rule-based tagger. However his system too used a lot of sequential information and did equally well (96.5% accuracy). Ratnaparkhi (1996) inserted a bunch of extra features for rarely occurring

and ambiguous words. To avoid feature space explosion, he used a maximum entropy Markov model. This was an ad-hoc strategy meant specifically to improve performance on the POS tagging problem. Recently Marquez & Padro (2000) presented a sliding window approach with an output context of 5. They divided the POS tagging problem into many classification problems, one for each ambiguity class and induced statistical decision trees to model grammar constraints (accuracy 99%). Lafferty, McCallum & Pereira (2001) employed the POS tagging domain to present Conditional Random Fields. POS tagging has applications in language understanding and information extraction (a problem in itself for which state-of-the-art systems learn sequential patterns of text).

2.4 Information Extraction

The field of information extraction (IE) deals with segmenting and extracting information from documents. For example, we might wish to extract information from web pages such as names, address, and title of the CEO of a company (extracted from a corporate webpage). Another example would be to extract the authors, title, journal, year, and page numbers of a published article by analyzing the citations in a paper. Since it is more likely that the name of the author would follow a word like "Author" or a phrase like "Written by" in the document, there is a possibility of extracting information from the sequence in which words appear. IE can therefore be modeled as a sequential supervised learning problem where the x 's represent the words in the document and the y 's represent a binary class (e.g., "author" and "not author"). Classical IE systems applied NLP as their core. Statistical methods were employed in IE systems like FASTUS 1993, and CIRCUS 1993. Cohen (1995) presented text classification using relational learning with ILP methods and Freitag (1996) proposed the use of Machine Learning methods with IE where text formats are fairly informal along with the use of sliding windows and HMMs for learning sequential text patterns. Later HMMs started coming into use more frequently in IE (Loek 1997) and Craven et al (1998) suggested the use of relational learning methods that learn to classify Web pages on the basis of the sub-graph of pages surrounding a given page. This was an application of IE where they identified useful fields in existing text (in their case, web pages). Freitag & McCallum (1999) implemented HMMs to model the same for IE with a statistical technique called shrinkage which, they explain, brought in the best of both worlds in terms of the bias variance tradeoffs. More recently McCallum, Freitag & Pereira (2000) exhibited the use of maximum entropy Markov models for IE. This technique not only modeled the causal relationships in the problem better but also reduced the space of possible feature value combinations required to be enlisted. Sequential supervised learning techniques are commonplace and taken for granted in IE applications today.

2.5 Summary

Among other applications of sequential data classifiers are handwriting recognition systems (Bengio et al. 1995, Hu & Brown 1996), Speech Recognition (Rabiner 1989), intrusion detection systems (Lee, Stolfo & Mok 1998), fraud detection systems (Fawcett & Provost 1997) and many more. Solutions to all these problems have been either ad-hoc or extremely complex. Since we can classify all the above problems as sequential supervised learning problems, we can imagine a generalized approach to solving them without using ad-hoc strategies. A general purpose tool is needed for such problems. Such a robust and efficient off-the-shelf implementation will not only aid in the development of new applications in these and other areas but also allow researchers to share a common platform for comparative study.

RSW is the first general purpose off-the-shelf classifier for sequential data using recurrent sliding windows. However, since a general purpose tool is, by definition, not tuned to solve any specific problem, the two issues discussed in section 1.4 arise. Both these issues can be summed up into one question: given a problem domain, how can we calibrate the window size for a recurrent sliding window classifier so as to achieve optimal performance? In this paper along with RSW we also present the first systematic study of calibration of window size for RSW through our experimental study.

3 Experiments and Results

3.1 Experimental Methods:

For our experiments, we chose two large datasets. The first is the text-to-speech dataset from the NETtalk system (Sejnowski & Rosenberg 1987). This dataset consists of 20,003 words, expressed as sequences of letters

and their corresponding pronunciations in terms of phonemic sound representations. Of these, 2000 sequences (words) were picked at random without replacement. This dataset of 2000 sequences was then randomly split into a 1000-sequence training dataset and a 1000-sequence test dataset. Each class label consists of a phoneme-stress pair. For example, in words like “there”, the “t” is pronounced as “T>” wherein the “T” is the phoneme label and the “>” is the stress label. There are 56 phonemic sounds produced in English speech and 6 different levels of stress. Of all the 336 possible combinations of phonemes and stresses only 140 are observed in the training and test sets. The average length of English words in the data is 7 (minimum 2, maximum 17).

The second dataset we studied was the Protein Secondary Structure Prediction dataset employed by Qian & Sejnowski (Qian & Sejnowski 1988). Their task was to assign each element in a protein sequence (drawn from a 20-letter alphabet) to one of three classes (alpha helix, beta sheet, or coil). Although there are only three classes, the average length of the protein sequences is 169 (minimum 11, maximum 498).

We computed two measures of performance:

- (a) The percentage of individual elements in the test set predicted correctly and
- (b) the percentage of whole sequences in the test set predicted correctly. For the latter, all elements in the sequence must be correctly predicted in order for the sequence to be declared correctly predicted.

Because the sequences are so long, we never observed a case where an entire protein sequence was correctly predicted.

For each dataset and each performance measure, we applied RSW with $LIC = RIC = 0 \dots 7$. This resulted in an input window size ranging from 1 to 15 by steps of 2. For each input window size, experiments were performed with the right output context (ROC) varying from 0 to 7 and then with the left output context (LOC) varying from 0 to 7. This gives a total of 120 experiments with different input and output context combinations.

Each of the 120 experiments was performed with each of four learning algorithms: naive Bayes, J48¹ decision trees, bagged J48 trees, adaboosted J48 trees (boosting by weighting). For decision trees, we employed pessimistic pruning with five confidence factors (0.1, 0.25, 0.5, 0.75, 1). With bagged and boosted trees, we experimented with 10 to 50 unpruned trees.

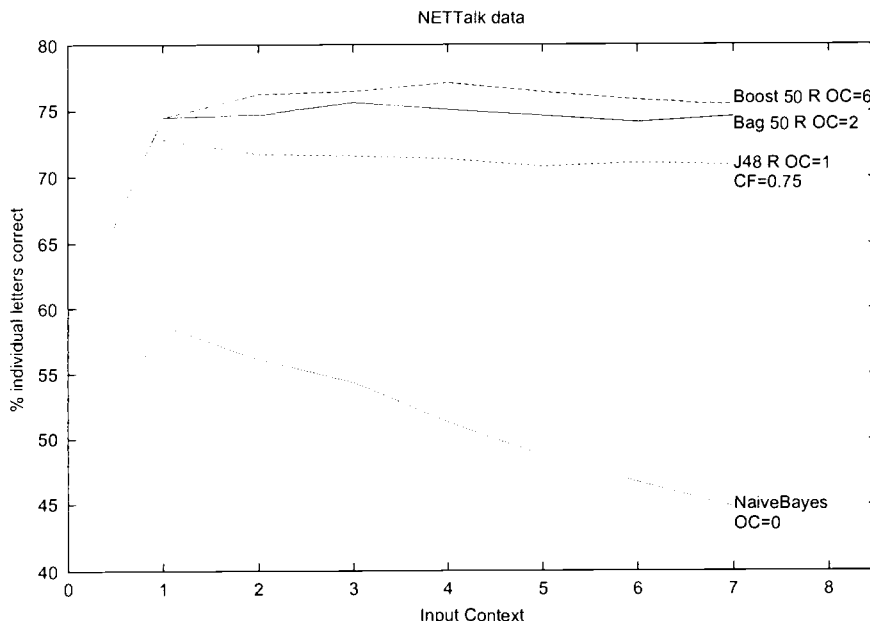


Figure 5: NETTalk: % of individual letters correct

¹ J48 is the WEKA implementation of C4.5 (Quinlan 1993)

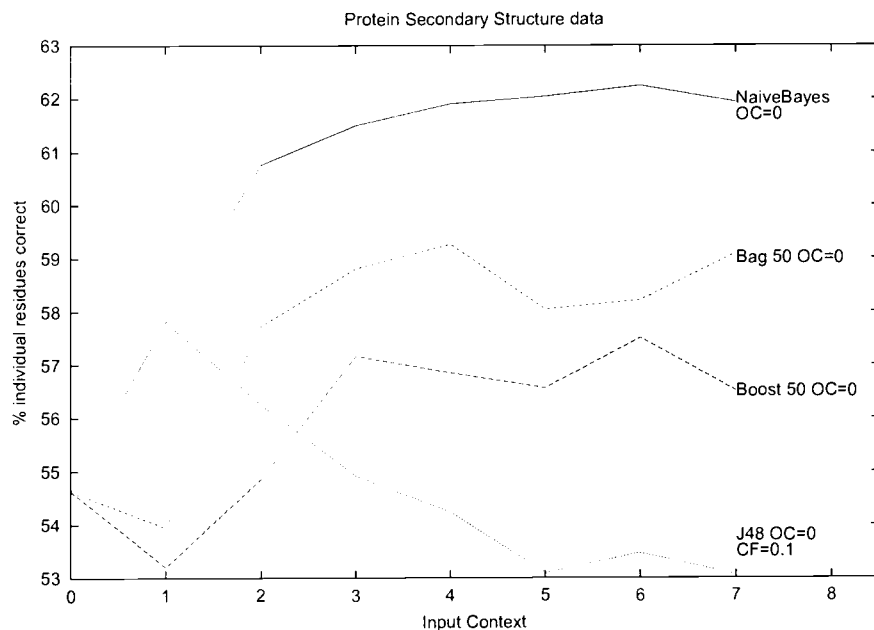


Figure 6: Protein Secondary Structure: % of individual residues (amino acids) correct

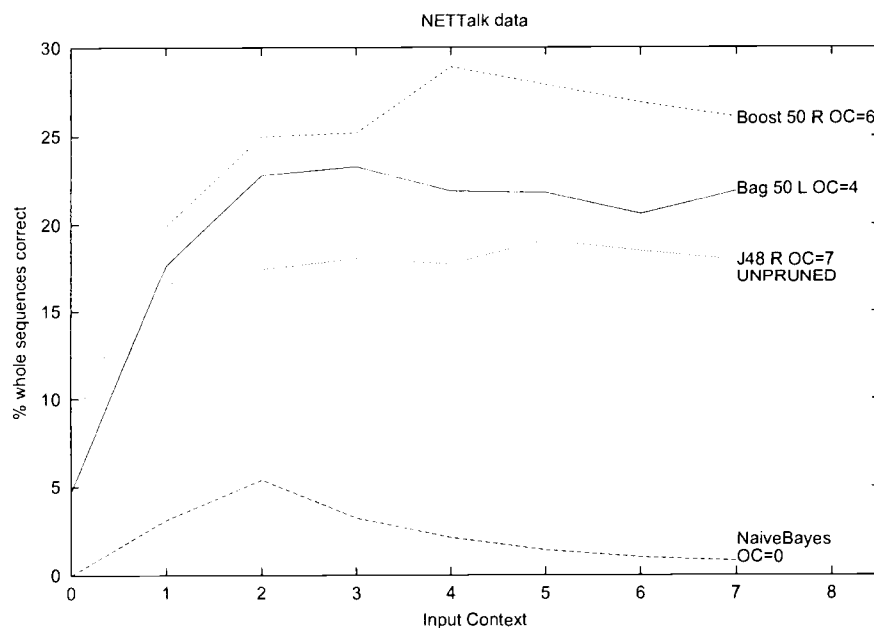


Figure 7: NETTalk: % of whole words correct

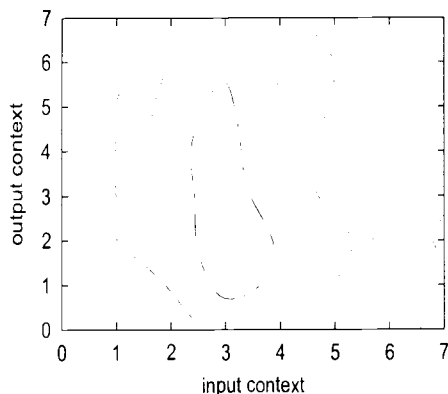


Figure 8: Contour plot showing Bagging performance as a function of input and output context.

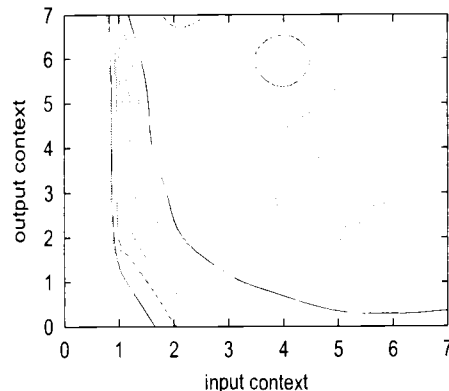


Figure 9: Contour plot showing boosting performance as a function of input and output context.

3.2 Results

Figures 5, 6, and 7 compile the results of the above experiments. Each graph displays a comparison of the performance of various learning algorithms across various input contexts for the best configuration of the remaining parameters. A label next to each performance curve displays the name of the learning algorithm and the output context. The output context chosen for each curve is the one where the performance of the algorithm achieves its maximum. Where there is a difference in performance between the use of left and right output context, it is specified appropriately by an *L* or an *R*.

Figure 5 shows performance curves for the four algorithms applied to the NETtalk data. The performance criterion is the number of letters in the test dataset predicted correctly. Naive Bayes works very poorly on this problem. Peak performance is obtained with no output context and an input context of 2 which is a simple 5-letter window. A clear case of overfitting is exhibited as the input context increases.

The best configuration for a single decision tree is the following: input context of 1 (i.e., 3-letter window), output context of 1, and a pruning confidence factor of 0.75. Performance starts dropping slightly for greater input contexts. Bagging achieves significantly better performance, and it attains maximum performance with an input context of 3, output context of 2 (and 50 bagged trees). Boosted decision trees do somewhat better than Bagging. They achieve maximum performance with an input context of 4 and an output context of 6. So there is a pattern: as we move from a single tree ($IC = 1, ROC = 1$) to Bagging ($IC = 2, ROC = 2$), to Boosting ($IC = 4, ROC = 6$), the learning algorithms are able to handle larger input and larger output context.

From this figure, we draw the following conclusions. First, the best configuration depends on the learning algorithm. Second, as expected, all of the algorithms exhibit some overfitting when the input context becomes too large. Third, all of the algorithms give very poor performance with an input context of 0 (i.e., only the one current letter in the window) which is equivalent to treating this problem as a standard supervised learning problem. Fourth, all algorithms obtain better results using right output context rather than left output context. This is consistent with the results of Bakiri (1991).

Figure 6 shows the performance of the four algorithms applied to the Protein dataset. In this domain, naive Bayes gives the best performance ($IC = 6$). The second best is bagged decision trees ($IC=4$). The third best is adaboosted decision trees ($IC = 6$), and the worst is a single decision tree with strong pruning ($IC = 1$, pruning confidence 0.1). Note that in all cases, an output context of 0 is preferred. Most state-of-the-art secondary structure prediction methods use non-recurrent sliding windows coupled with better input features and additional post-processing of the predictions (Jones, 1999).

Figure 7 shows the four algorithms applied to the NETtalk data again, but this time the evaluation criterion is the number of whole words in the test set predicted correctly. This graph is somewhat similar to the one in Figure 5, except that the peak performances now have different window sizes. Bagging does better with a smaller input context and a larger output context. Boosting uses the same output context of 6 but performs better with an input context of 4 instead of 3. A big surprise is that J48 gives its best performance with an output context

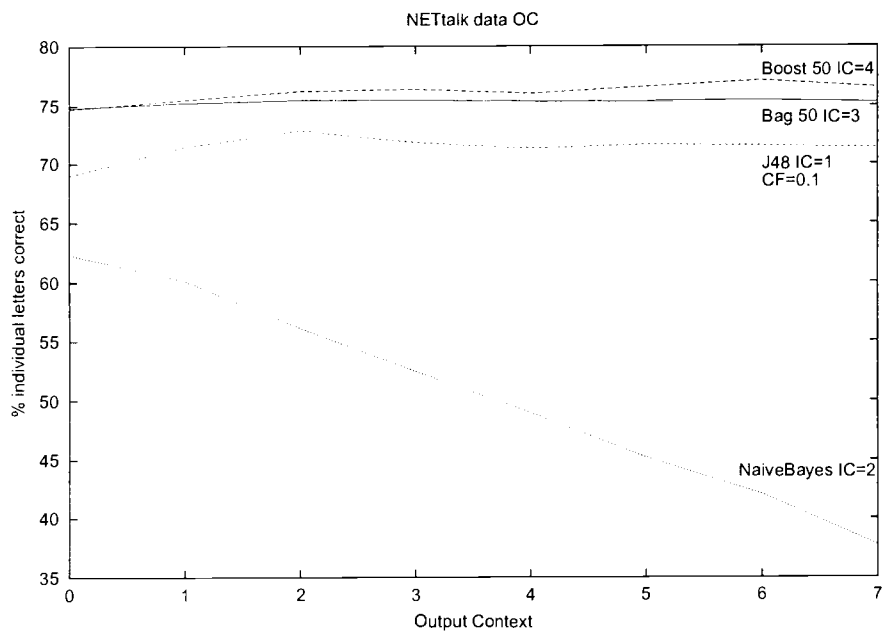


Figure 10: NETTalk: % of individual letters correct as a function of output context

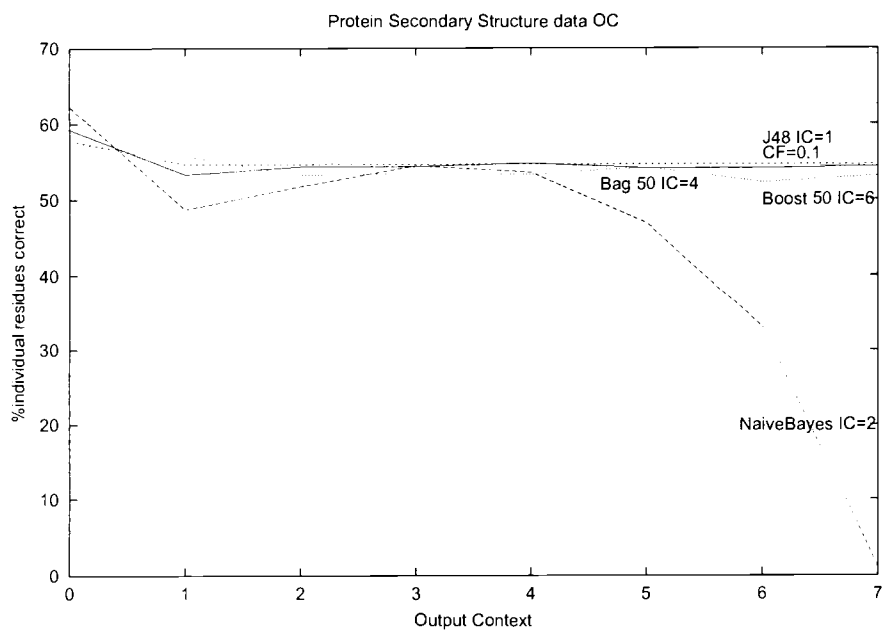


Figure 11: Protein Secondary Structure: % of individual residues (amino acids) correct as a function of output context

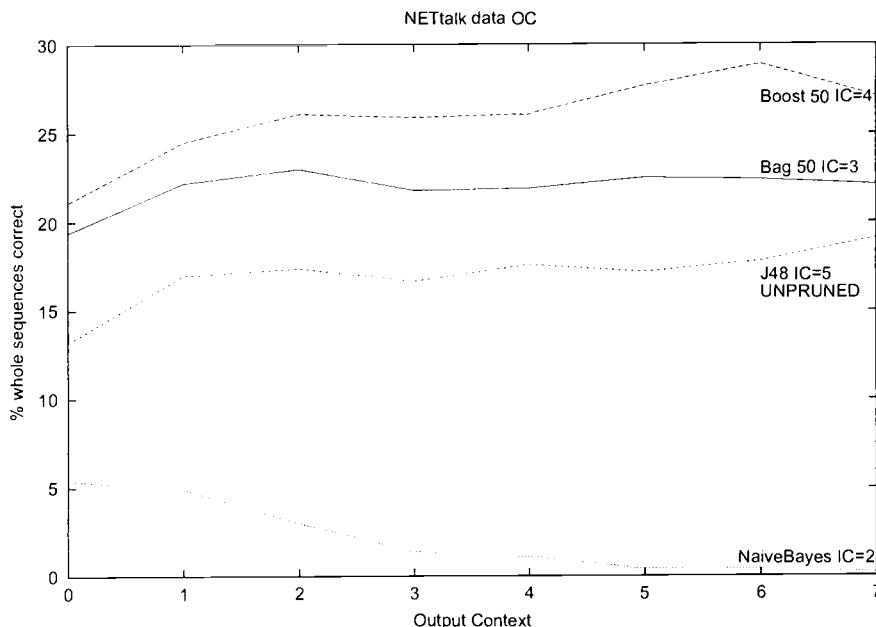


Figure 12: NETTalk: % of whole words correct as a function of output context

of 7 (and no pruning). Even naive Bayes changes its best configuration from input context 1 to input context 2. The major conclusion to be drawn from this figure is that the optimal choice of input and output context depends on both the learning algorithm and the evaluation criterion.

Figures 8 and 9 show contour plots of the performance of Bagging and Boosting for various input and right output contexts. For Bagging, the maximum at input context 3, output context 2 is clearly visible as is the maximum at input context 4 and output context 6 for Boosting. There is some evidence for a trade-off between input and output context. For example, decreasing input context does allow us to increase output context or vice versa to a certain extent. This is the case, in Figure 9. This is less evident in Figure 8. However, there is a large plateau of good performance when the input and output contexts are sufficiently large. Here shrinking either input or output context hurts performance.

Figures 10 and 11 show the performance varying with output context for the different learning algorithms. The performance criterion is the number of individual examples predicted correctly. Here again all other parameters including the input context are chosen so as to maximize the performance. The input context for every curve is specified in a label adjacent to it.

We observe here that the curves for Bagging, Boosting and J48 are nearly flat, which shows that the amount of output context has little impact on their performance. In comparison, Figure 5 showed that the amount of input context has a much greater impact in the sense that it must be at least 1 in order to get a reasonable performance. However, for naive Bayes, increasing the output context has a devastating effect — performance drops rapidly and monotonically as the output context is enlarged.

The picture changes when we consider the percentage of whole sequences correctly classified as seen in Figure 12. Here, increasing the output context improves the accuracy of J48 and Boosting substantially, and it improves Bagging to a lesser extent. Again the accuracy of naive Bayes decreases with increasing output context. Clearly naive Bayes overfits for large input and output window sizes.

4 Discussion

It is clear from the results obtained that the optimal window size depends on a lot more than the underlying domain. In this section we attempt to explain why this is so. The problem domain certainly has information

to provide, but it is not sufficient for calibrating a sequential data classifier. In order to explain the results we turned to the popular and widely accepted bias-variance theory. Bias-variance theory employs the standard statistical analysis that views each training set as a random sample from a population of possible training sets. It decomposes the prediction error of a classifier into three factors: (a) the bias, or systematic error of the learning algorithm, (b) the variance, or variation of the learning algorithm that results from random sampling of the given data set from the population, and (c) the noise, or variation in the class labels due to factors that cannot be predicted from the input features. In continuous prediction problems where the error is measured by the squared difference between the predicted and observed values of the response variable, the error can be decomposed into a sum of the squared bias, the variance, and the noise. In classification problems where the error is measured by the error rate, the decomposition is more complicated, and there has been considerable debate about the most appropriate bias-variance-noise decomposition. We will employ the decomposition developed by Domingos (2000). He developed his decomposition for the case where there are only two possible classes. We will first present his decomposition and then show how to extend it for more than two classes. Consider that you have lots of different

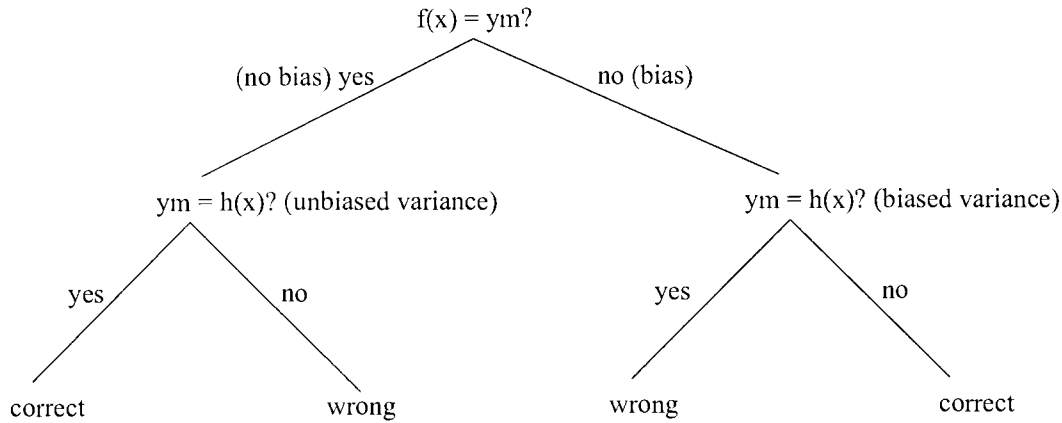


Figure 13: A Case Analysis of Error for 2 classes

training sets (of equal size) drawn from the same population at random. Consider that you have a test set that is representative of the population. Consider a classifier h trained on each of the training sets and tested on each example (x, y) of the test set. Let y_m be the majority prediction among all the classifiers for (x, y) . Also consider that the true function we are trying to learn is f . There are 3 components that determine whether the prediction $y = h(x)$:

Noise: $y = f(x)$? (Is the label on x correct?)

Bias: $f(x) = y_m$? (Is the majority prediction among all classifiers correct?)

Variance: $y_m = h(x)$? (Does the majority prediction y_m match the prediction of the particular hypothesis h ?)

To simulate this analysis using real datasets, we generate 200 bootstrap replicates of the training data (X, Y) . On each bootstrap replicate, we train a classifier h using a learning algorithm that attempts to learn the true function f . Then, we classify each test point x on each classifier. We assume that the noise is zero.

Domingos developed a case analysis of error. If there is no noise, this analysis can be shown as in Figure 13. The first level in this tree structure checks for bias, and the second level checks for variance. In the case where the classification is not biased, if there is no variance, then the prediction is correct. Otherwise variance causes error. This type of variance is called unbiased variance and it is undesirable. In the case where the classification is biased, lack of variance is bad (biased variance in this case), because biased variance will cause the prediction to differ from the biased classification. In the 2-class case, if the prediction is not equal to the wrong class, it has to be equal to the right class. Therefore biased variance is desirable variance. Over the entire test set, the total bias, unbiased and biased variance is calculated by simple summation and the error is measured as follows:

$$Error(l) = \sum Bias(b) - \sum biasedVariance(V_b) + \sum unbiasedVariance(V_u)$$

In order to perform a bias variance analysis of our results, we extended Domingos' work for the case where there are multiple classes. A case analysis of error for the multiple class case can be shown as in Figure 14. Once

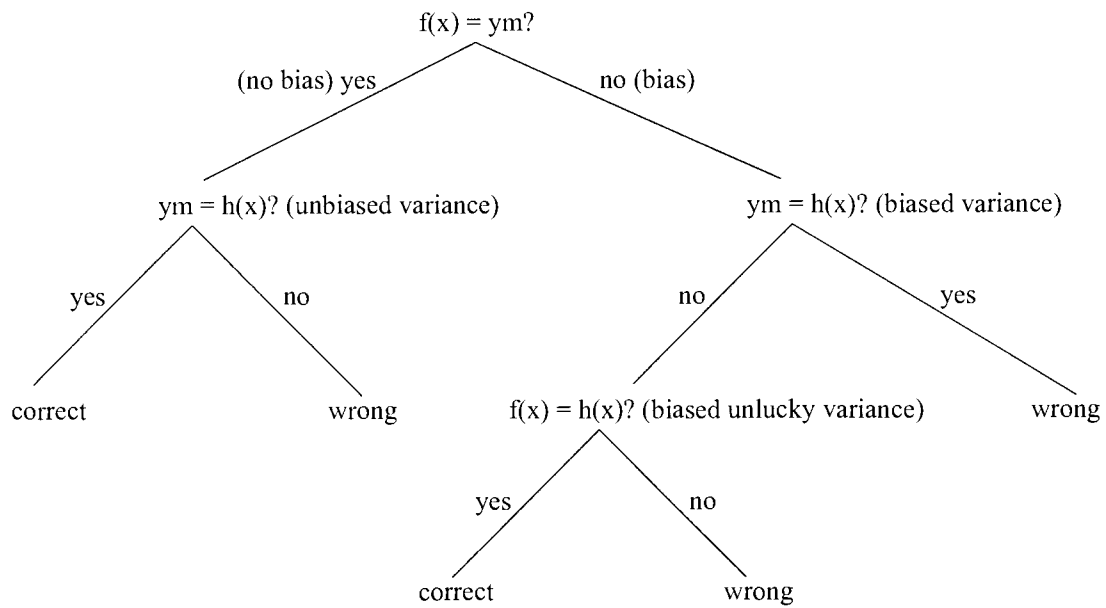


Figure 14: A Case Analysis of Error for multiple classes

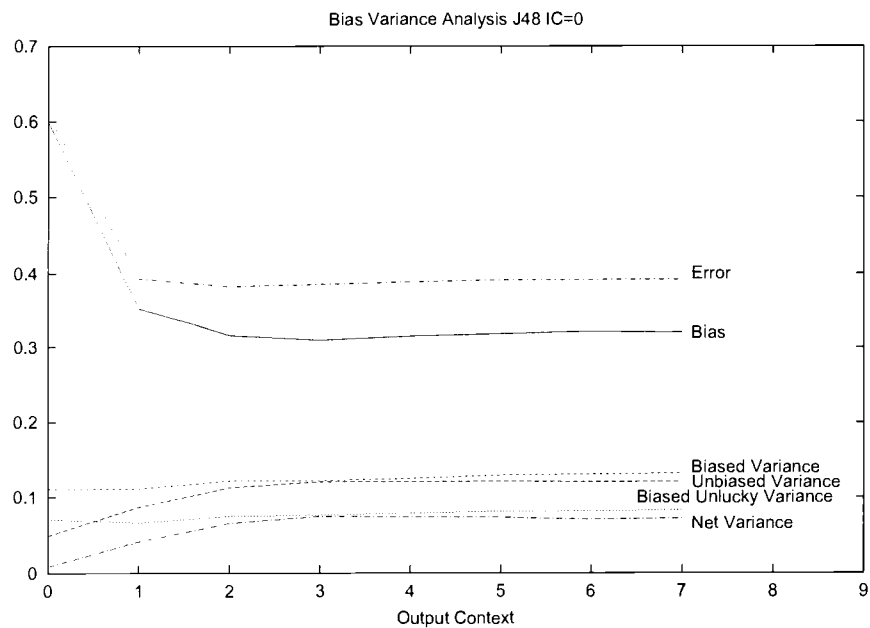


Figure 15: Bias Variance Analysis of J48 against output window size with IC=0 for NETtalk data

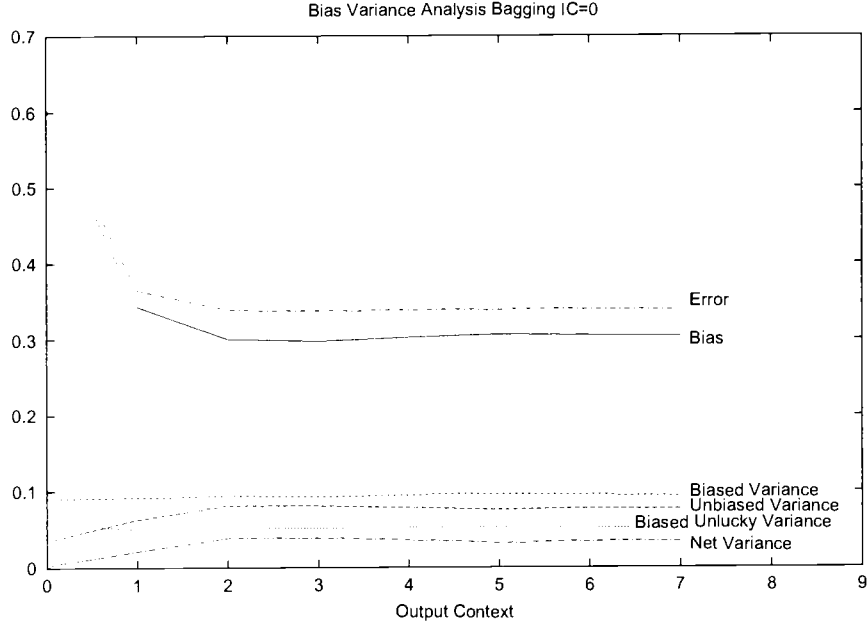


Figure 16: Bias Variance Analysis of Bagged J48 against output window size with IC=0 for NETtalk data

again we assume there is no noise. The only difference from the 2-class setting is for the biased variance case. In this case, the classification is biased and suffers from variance as well. But this biased variance is not necessarily good, because with more than 2 classes now, it is not enough for the prediction to be unequal to a wrong answer. Since there are multiple wrong answers here, the prediction can be biased, differ from the majority prediction, and still be wrong. This brings us to another kind of variance we term “biased unlucky variance”. It is unlucky, because even when the classification decision facing this kind of variance differs from the incorrect majority, it is wrong and hence loses on both fronts. The error decomposition now has a new term and can be written as follows:

$$Error(l) = \sum Bias(b) - \sum biasedVariance(V_b) + \sum unbiasedVariance(V_u) + \sum BiasedUnluckyVariance(V_{bu})$$

Figures 15 and 16 show the results of our experiments on the bias-variance analysis of single decision trees and 50 bagged trees respectively. For the purpose of these experiments, we have chosen the NETtalk data and assumed there is no noise in the data. The Figures 15 and 16 plot bias, variance, and error rate against output context for no input context. It can be seen from both these curves that bias decreases with increasing context (window size) and net variance increases with increasing context (window size). This corresponds to our intuition. Suppose the data was generated by a second order hidden Markov model. If a learnt classifier now used an output context of 0 or 1, there will be a bias, because the learnt classifier cannot represent the true function. The greater the window size, greater is the probability that the true order or the generative model is encompassed by the learnt classifier. Hence there is a lower bias. At the same time, a large window equates to a large number of features in the sliding window realm. Variance increases with the number of features in the data set and hence should also increase with increasing window sizes.

Figures 17, 18, and 19 compare the error rate, bias, and net variance curves of J48 decision trees and 50 bagged trees for respectively 3 different output contexts 0, 4, and 7. Each of the plots tell the same story. Bagging J48 decision trees gives a reduction in net variance as it should according to theory. The apparent small reduction in bias with Bagging is actually because of the fact that to optimize performance the bagged trees have been left unpruned whereas single decision trees have been pruned using pessimistic pruning with a confidence factor of 0.25. Pruning reduces variance but increases bias in trees. With Bagging however, the aggregation reduces variance, so better performance is obtained by growing the trees to full depth and thereby reducing bias. The bias curves for Bagging and J48 are almost overlapping indicating there is not much difference in the bias for

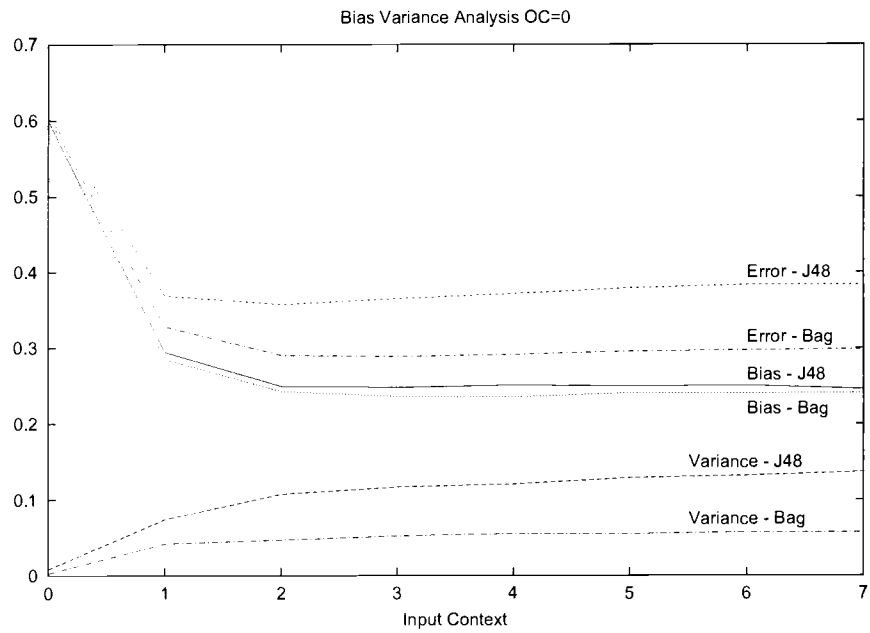


Figure 17: Bias Variance Analysis of J48 and Bagging as a function of input window size with OC=0 for NETtalk data

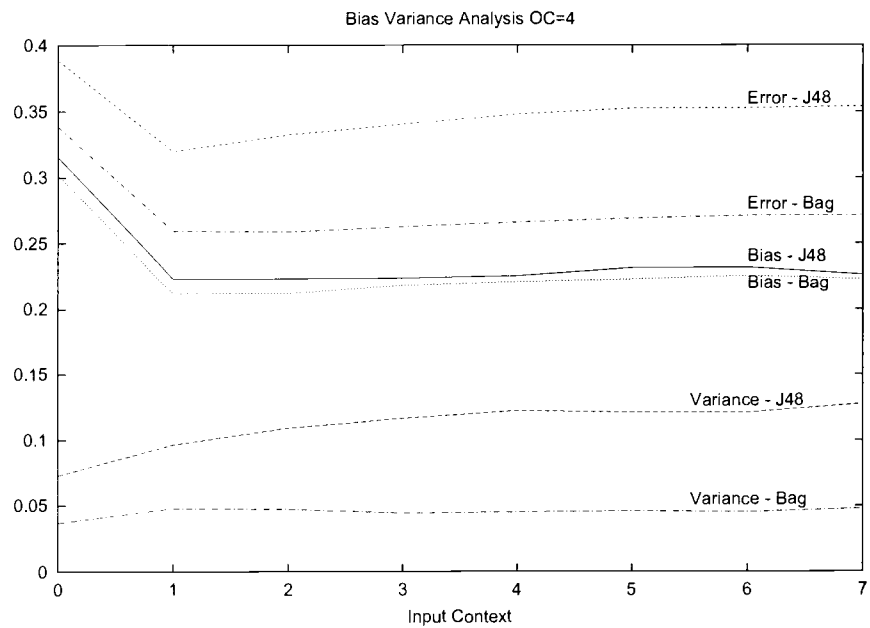


Figure 18: Bias Variance Analysis of J48 and Bagging as a function of input window size with OC=5 for NETtalk data

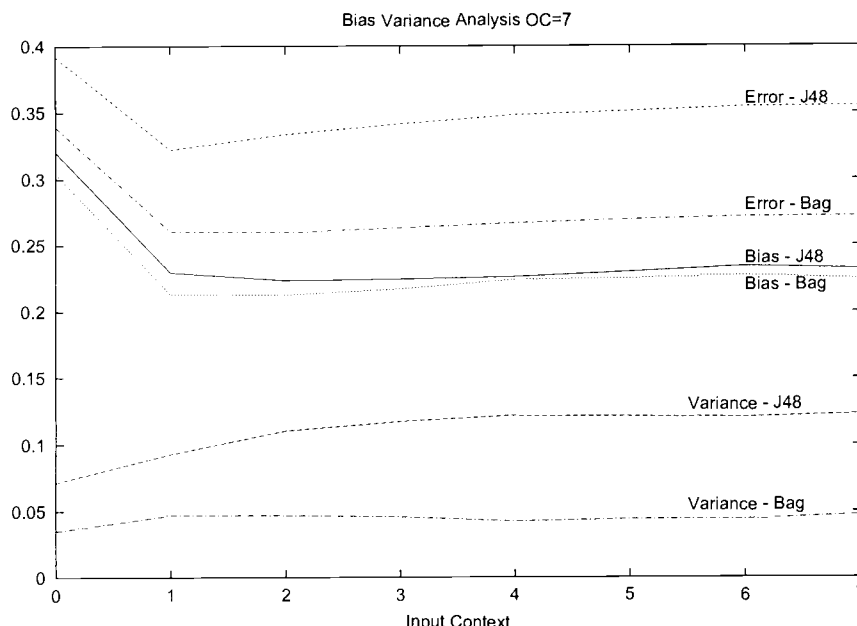


Figure 19: Bias Variance Analysis of J48 and Bagging as a function of input window size with OC=7 for NETtalk data

all window sizes. Bias decreases steadily as the window size increases. It decreases up to an input context of 3 and then remains more or less constant. The variance curves are a lot different. There is a huge gap in variance between J48 and Bagging, and this difference is increasing with increasing window size. Clearly Bagging is able to eliminate the variance caused by using single decision trees. Variance due to bagged decision trees increases much slower than variance due to J48. The optimal operating point (the point where error is minimum) is at an input context of 2 for J48 and an input context of 3 for Bagging in Figure 17. Bagging, therefore, employs a larger window and obtains a better performance. Similarly in Figures 18 and 19, the window size at which Bagging reaches its optimal point is greater than the window size at which J48 reaches its optimal point.

This is the point where the story gets interesting. With standard supervised learning, sources of bias and variance for the classifier are the training data and the learning algorithm which generates the classifier. However with sequential supervised learning, as we saw earlier, the window size is an additional source of bias and variance. A bagged classifier can afford a larger window size, because Bagging has lower variance. The extra variance generated by a greater window size is eliminated by bagging the classifier. The bagged classifier then uses this larger window size to obtain some additional bias reduction. This is why Bagging can use a larger window size and yet achieve better performance. The same argument can be extended to Boosting, which is a bias and variance reduction technique. This confirms the idea that differences in optimal window sizes for different learning algorithms are due to the different bias variance profiles of these algorithms.

5 Conclusions and Future work

All methods for solving sequential supervised learning problems must confront the question of how much context to employ. For HMMs and conditional random fields, the question of output context becomes the question of the order of the Markov process (first order, second order, etc.). HMMs cannot handle large input contexts (at least not without violating the assumptions of the generative model), but CRFs can. Hence, CRFs also must face the choice of the size of the input context. The choice of input and output context is analogous to classical feature selection. One common approach to feature selection is to fit some simple model to the data and/or compute some figure of merit for the informativeness of each attribute (Kira & Rendell, 1992; Koller & Sahami, 1996).

The experiments reported here show that this will not work for SSL problems, because the correct choice of input and output contexts depends on the learning algorithm. Although a simple bias-variance analysis explains the results, the input and output contexts chosen by some simple method (e.g., naive Bayes) are not good choices for boosted decision trees. A consequence of this is that the input and output contexts need to be chosen by cross-validation or holdout methods.

The experiments have also shown that the optimal input and output contexts depend on the performance criterion. Window sizes that maximize the correct classifications of individual examples are not the ones that maximize the correct classifications of whole sequences and vice versa.

An important question for future research is to understand why different window sizes are required for different performance criteria (whole sequences versus individual items). Our bias-variance analysis explains why Bagging and Boosting are able to handle large window sizes, but it does not explain why larger output contexts are more valuable for whole sequence classification. Intuitively, whole sequence classification requires greater integration of information along the sequence. The primary way of integrating information along the sequence is through the output context. But this intuition needs to be made more precise in order to obtain a complete theory of sequential supervised learning.

Acknowledgements

The authors gratefully acknowledge the support of the National Science Foundation under grants IIS-0083292 and ITR-0001197.

References

- Bakiri, G. (1991). Converting English text to speech: A machine learning approach. Tech. rep. 91-30-2, Department of Computer Science, Oregon State University, Corvallis, OR.
- Craven, M., DiPasquo, D., Freitag, D., McCallum, A., Mitchell, T., Nigam, K., & Slattery, S. (1998). Learning to extract symbolic knowledge from the World Wide Web. In *Proceedings of the 15th National Conference on Artificial Intelligence (AAAI-98)*, pp. 509–516 Menlo Park, CA. AAAI Press/MIT Press.
- Jones, D. T. (1999). Protein secondary structure prediction based on position-specific scoring matrices. *Journal of Molecular Biology*, 120, 97–120.
- Kira, K., & Rendell, L. A. (1992). A practical approach to feature selection. In *Proceedings of the Ninth International Conference on Machine Learning*, pp. 249–256 San Francisco, CA. Morgan Kaufman.
- Koller, D., & Sahami, M. (1996). Toward optimal feature selection. In *International Conference on Machine Learning*, pp. 284–292.
- Lafferty, J., McCallum, A., & Pereira, F. (2001). Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the 18th International Conference on Machine Learning*, pp. 282–289 San Francisco, CA. Morgan Kaufmann.
- Màrquez, L., Padró, L., & Rodríguez, H. (2000). A machine learning approach to POS tagging. *Machine Learning*, 39(1), 59–91.
- Qian, N., & Sejnowski, T. J. (1988). Predicting the secondary structure of globular proteins using neural network models. *Journal of Molecular Biology*, 202, 865–884.
- Quinlan, J. R. (1993). C4.5: Programs for Empirical Learning. Morgan Kaufmann, San Francisco, CA.
- Rabiner, L. R. (1989). A tutorial on hidden Markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2), 257–286.
- Sejnowski, T. J., & Rosenberg, C. R. (1987). Parallel networks that learn to pronounce english text. *Journal of Complex Systems*, 1(1), 145–168.

- McCallum, A., Freitag, D. & Pereira, F. (2000). Maximum Entropy Markov Models for Information Extraction and Segmentation. In *Proceedings of the 17th International Conference on Machine Learning*, pp. 591–598 Stanford, CA. Morgan Kaufmann.
- Brill, E. (1994). Some Advances in Transformation-Based Part of Speech Tagging. In *Proceedings of the 12th National Conference on Artificial Intelligence*, pp. 722–727 Seattle, Washington. AAAI Press/MIT Press
- Cutting, D., Kupiec, J., Pedersen, J., & Sibun, P. (1992). A Practical Part of Speech Tagger . In *Proceedings of the 3rd Conference on Applied Natural Language Processing*, pp. 133–140 Xerox PARC technical report SSL-92-01
- McCallum, A., Freitag, D. (1999). Information Extraction with HMMs and shrinkage . In *Proceedings of the AAAI-99 Workshop on Machine Learning for Informatino Extraction*, pp. 31–36 Orlando, FL
- Bengio, Y. (1999). Markovian Models for Sequential Data. In *Neural Computing Surveys*, 2, pp. 129–162
- Ratnaparkhi, A. (1996). A Maximum Entropy Model for Part-Of-Speech Tagging. In *Proceedings of Conference on Empirical Methods in Natural Language Processing*, pp. 133–141 Philadelphia, 1996
- Witten, I. H., & Frank, E. (1999). *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*. Morgan Kaufmann, San Francisco, CA.